



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2012

Revealing the complexity of health determinants in resource poor settings

Lewis, F I ; McCormick, B J J

Abstract: An epidemiological systems analysis of diarrhea in children in Pakistan is presented. Applying additive Bayesian network (ABN) modeling to data from the Pakistan Social and Living Standards Measurement (PSLSM) survey reveals the complexity of child diarrhea as a disease system. The key distinction between standard analytical approaches, such as multivariable regression, and Bayesian network analyzes is that the latter attempts not only to identify statistically associated variables, but to additionally, and empirically, separate these into those directly and indirectly dependent with the outcome variable. Such discrimination is vastly more ambitious but has the potential to reveal far more about key features of complex disease systems. Additive Bayesian network analyzes across 41 variables from the PSLSM identified 182 direct dependencies, but with only three variables: Access to a dry pit latrine (protective: OR=0.67); Access to an atypical water source (protective: OR=0.49); and No formal garbage collection (unprotective: OR=1.32), supported as directly dependent with the presence of diarrhea. All but two of the remaining variables were also in turn directly or indirectly dependent with these three key variables. These results are contrasted with the use of a standard approach (multivariable regression).

DOI: <https://doi.org/10.1093/aje/KWS183>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-62970>

Journal Article

Accepted Version

Originally published at:

Lewis, F I; McCormick, B J J (2012). Revealing the complexity of health determinants in resource poor settings. *American Journal of Epidemiology*, 176(11):1051-1059.

DOI: <https://doi.org/10.1093/aje/KWS183>

Revealing the complexity of health determinants in resource poor settings

Fraser I. Lewis*, University of Zurich, Section of Epidemiology, Winterthurerstrasse 270, 8057 Zurich, Switzerland, Phone +41 44 635 9051. Email: fraseriain.lewis@uzh.ch

*corresponding author

Benjamin J. J. McCormick, Fogarty International Center, National Institutes of Health, 31 Center Drive, Bethesda, Maryland 20892, USA, Phone + 01 301 496 0815. Email : ben.mccormick@nih.gov

March 19, 2012

Abstract

An epidemiological systems analysis of diarrhea in children in Pakistan is presented. Applying additive Bayesian network (ABN) modeling to data from the Pakistan Social and Living Standards Measurement (PSLSM) survey reveals the complexity of child diarrhea as a disease system. The key distinction between standard analytical approaches, such as multivariable regression, and Bayesian network analyzes is that the latter attempts not only to identify statistically associated variables, but to additionally, and empirically, separate these into those directly and indirectly dependent with the outcome variable. Such discrimination is vastly more ambitious but has the potential to reveal far more about key features of complex disease systems. Additive Bayesian network analyzes across 41 variables from the PSLSM identified 182 direct dependencies, but with only three variables: Access to a dry pit latrine (protective: OR=0.67); Access to an atypical water source (protective: OR=0.49); and No formal garbage collection (unprotective: OR=1.32), supported as directly dependent with the presence of diarrhea. All but two of the remaining variables were also in turn directly or indirectly dependent with these three key variables. These results are contrasted with the use of a standard approach (multivariable regression).

Key words: Epidemiologic Determinants; Diarrhea; Bayesian Network;
Graphical Model; Socioeconomic Factors

Presented here is an epidemiological systems approach for identifying potential determinants of diarrhea in children under five years using data from the Pakistan Social and Living Standards Measurement (PSLSM) survey. Childhood diarrhea is the second biggest cause of worldwide mortality in children under five years (1, 2) and health surveys targeting this disease are common (3–6). While such study designs are not without issue, potentially suffering from data quality and reliability concerns (e.g.(7)), they are widely used as low cost methods of data collection in developing countries.

A major challenge when analyzing data from surveys is that they are typically exploratory in nature, where the precise aetiology of health outcomes are not known, and information on a large number of variables is often collected, not all of which are necessarily complete or relevant. It is also typical that many variables of potential interest are inter-related, both to each other and also to health outcomes. Such data may be conceptualized as describing an epidemiological system (8–11), that is, a collection of mutually inter-dependent variables some or all of which can predict or affect the health outcomes of interest.

Additive Bayesian networks (ABN) are introduced as a methodology for identifying statistical dependencies in complex disease systems using observational data. Ultimately, what is desired in many epidemiological analyzes is the identification of causal pathways (e.g. (12), (13, 14)), which can be extremely challenging in systems such as diarrheal disease where many high level casual factors have been postulated(15), and the identification of statistical dependencies can be invaluable for informing such analyzes.

The key distinction between standard multivariable regression analyzes and

BN type analyzes is that multivariable regressions seek to identify covariates associated with some outcome variable, e.g. presence of diarrhea. Bayesian network analyzes go much further and attempt not only to identify associated variables, but to additionally, and crucially empirically, separate these into those directly and indirectly dependent with the outcome variable. The latter is vastly more ambitious but has the potential to reveal far more about key features of complex disease systems than existing commonly used approaches. This is the central message of the work presented: ABN analyzes are superior to standard approaches for inferring statistical dependencies from complex observational data.

Identifying statistical dependencies using multivariable regression

When performing exploratory analyzes of data comprising of many variables, it is common to utilize some form of multivariable regression in which a variable selection process is then employed. The goal being to search for variables which are statistically significantly associated with, say, an outcome variable such as disease presence. Stepwise regression searches are widely used (e.g (16–19)) despite being viewed rather negatively in the epidemiological and biostatistical literature (20–22). Such automated searches are arguably over-used, or rather, that the results from such analyzes are too often presented without sufficient additional checks to ensure the robustness of associations against overfitting (23).

In rapidly developing and increasingly data rich fields such as genetic epidemiology, computational biology and bioinformatics, automation in statistical modeling is standard, and indeed arguably essential when faced with exploring observations from large numbers of potentially inter-dependent variables. That

automated searches tend to over-fit is well known. There are, however, well established techniques for addressing this; two of the most commonly utilized are model averaging (e.g. (24, 25)) and parametric bootstrapping (e.g. (26)), both of which are explored in the later case study analyzes.

Identifying statistical dependencies using Bayesian networks

Bayesian network analysis is a form of statistical modeling which derives, from empirical data, a graphical network describing the dependency structure between variables, where this is formally depicted as a directed acyclic graph (DAG). Bayesian networks are widely used in areas such as systems biology (27–29), in HIV and influenza research (30–33), and also analyzes of complex disease systems (34–37). The origins of BN modeling lie within the machine learning and data mining literature (27, 38) with an accessible non-technical introduction in (28).

In multivariable regression analyzes the goal is to identify statistically significant associations between an outcome variable and one or more covariates. Here “association” denotes that the variables are statistically dependent, and says nothing of whether the variables are directly or indirectly dependent. To borrow an example from (14), in multivariable regression analyzes with “lung cancer” as the outcome variable, and “smoking” and “yellow fingers” as covariates, then it may be expected that one or both of these covariates would be identified as statistically significantly associated with lung cancer. In contrast, in a Bayesian network analysis it would be expected that smoking and lung cancer be identified as directly dependent, and yellow fingers and smoking as directly dependent, but that yellow fingers and lung cancer not be identified as directly dependent. In

terms of a DAG this would describe a model with two arcs: one between lung cancer and smoking; and a second between yellow fingers and smoking - but with no arc between yellow fingers and lung cancer. Note we have *not* specified the direction of these arcs. In a BN analyzes each DAG is formally a factorization of the joint probability distribution of the observed data, and due to likelihood equivalence it is the presence of arcs between variables and not their direction which is the notable feature.

Consider the joint probability of variables X and Y , $P(X, Y)$. Theory gives $P(X, Y) = P(X|Y)P(Y)$ and $P(X, Y) = P(Y|X)P(X)$, where the former can be depicted as a DAG with one arc, from Y to X , and the latter with one arc from X to Y . The practical implication of this is that using observed data alone, it is not possible to statistically discriminate between different DAGs from within the same likelihood equivalence class - as these are probabilistically identical. However, determining likelihood equivalence between DAGs is extremely difficult in all but the simplest cases (see Web Appendix 1 for more details). Due to these complications it is typical to ignore arc direction in BN analyses (e.g. see (30–33, 39)), although notable exceptions are analyses of longitudinal data where dynamic Bayesian networks may be utilized (40). Using prior belief to impose explicit arc direction may be of some value in analyses which attempt to combine statistical dependency with causality and this is returned to later.

Material and Methods

Case Study Data

The Pakistan Social and Living Standards Measurement (PSLSM) survey is a

biennial survey of a large number of social, environmental and economic indicators, motivated and directed in part by efforts to meet the UN Millennium Development Goals. The survey is conducted at household level, sampling from the majority (approximately 97%) of the population across Pakistan. In the data analyzed here from 2005-2006 (Pakistan Social and Living Standards Measurement survey, 2004-2005. Federal Bureau of Statistics. Islamabad, Pakistan) 15,453 households were surveyed comprising 110,909 individuals of which 18,202 were under five years old. The survey includes around 250 questions, from which a broad subset (40) were included in the following analyzes based on potential relevance to childhood diarrheal disease determined from their inclusion in previous studies (see Web Appendices 2 and 3 for variable descriptions and details of previous studies). Diarrheal presence was taken from a binary question of whether children under five years old in the household had experienced diarrheal symptoms within the preceding 30 days.

Bayesian Network Modeling Formulation

Standard multivariable statistical models - linear or generalized linear models or their variants - are additive, that is, they describe the mean value of some response variable conditional on a given covariate pattern, as an additive contribution from each covariate. In contrast, BN models for categorical data, the most commonly utilized form of BN, use a parametrization where each and every covariate pattern is modeled using *independent* sets of parameters, that is, the parameters cannot be interpreted as main effects or interaction effects (see Web Appendix 4). This formulation may also be far from parsimonious(41) and does not provide any

ready interpretation of the model parameters. The standard formulation of BNs (e.g. (27)) does facilitate conjugacy, that is, all parameter estimates in a BN can be computed analytically for the three usual types of BN (categorical, Gaussian, and a special variant of mixed categorical and Gaussian variables(42)).

At the cost of a loss of conjugacy it is possible to formulate BN models which are direct analogues of standard multivariable linear and generalized linear models, where each variable in the data is modeled by an additive multivariable regression model, with an appropriate link function (e.g. a logit) if required (see Figure 1). As in classical BN models, additive BNs are described by a DAG. The price for this considerably greater model flexibility is that the goodness of fit and model parameters now require to be estimated numerically rather than analytically (Laplace approximations (43) are used here).

Bayesian inference requires prior distributions, and in a BN there are two possible types of priors: priors on the model parameters and priors on the DAG. In terms of parameter priors, the approach utilized here assumes uninformative Gaussian priors with zero mean and large variance for each of the regression parameters across all parts of the model, and diffuse Gamma priors for the precision parameter in Gaussian nodes of the model. In terms of structural priors it is currently assumed that each DAG structure is equally plausible in the absence of any data. Imposing prior causal knowledge onto network structures, e.g. by imposing conditions on arc direction, is returned to later.

Model selection in statistical analyzes

Statistical model selection comprises of three parts: A) choosing a general form

of model; B) deciding the scope of the model search space and how to cross it; and C) deciding how to summarize the results from B).

For the PSLSM case study, the first analysis presented comprises of a conventional stepwise regression search. Therefore, A) standard multivariable logistic regression; B) stepwise search, forwards from a null model and backwards from a full model with comparisons performed within a maximum likelihood framework, as that is what is provided in common statistical software and most usual in practice, and with AIC as the goodness of fit metric. For C) the single best model found in B) was then subjected to parametric bootstrapping to identify any issues of overfitting (see Web Appendix 5 for details).

In the BN analyzes, for A) an additive form of BN is used. For B) the two most widely utilized “structure discovery” approaches are used. Firstly, the local search heuristic detailed in (27) which is analogous to the usual stepwise search in multivariable regression. Secondly, a search over node orderings rather than DAG structures. Order based approaches were introduced in (44) and then substantially extended in (45). The motivation behind local heuristics (including stepwise searches in multivariable regression) is that they will identify high scoring - well fitting - models when it is not computationally feasible to identify the very best model with any certainty. The second approach for searching for optimal BN models is to collapse DAGs over node orderings; a node ordering is simply a list of the nodes, say as indices 1 through n , where a given DAG structure is consistent with an ordering if, and only if, the parents of each node precede their child node in this list. Orderings can be thought of as groups of DAG structures - those structures which are consistent with that particular ordering, and note

that each DAG may be consistent with more than one ordering, e.g. the empty DAG (no arcs) is consistent with every possible ordering. The basic idea is that by searching across orders the dimension of the search space is vastly reduced from $\approx n!2^{\binom{n}{2}}$ unique DAGs down to $n!$ unique orders(45), although the latter may still be computationally impractical. The price for this reduction in size of search space is that searching across orders is biased relative to searching across DAGs.

Finally, consider step C), how to summarize the results of BN model searches. Two options are either to construct some form of summary or “average” model by pooling across heuristic search results, or else select a single “best” model. A popular approach for the former is to construct a majority consensus network which builds a DAG comprising of all those arcs present in at least a majority ($> 50\%$) of the DAGs identified using heuristic searches(30, 35). Due to likelihood equivalence it is common to collapse over arc direction to avoid missing important structural features. For example, if arc $X \rightarrow Y$ appears in 50% of heuristic results, and $Y \rightarrow X$ appears in the other 50%, then even although this direct dependency between X and Y features in every search result it will never appear in a (directed) majority consensus network. For this reason collapsing over arc direction when presenting results of BN analyzes is common (e.g. (30, 39)). The purpose of summarizing over many DAGs is to address concerns of overfitting, and is directly analogous to the ubiquitous use of majority consensus trees in phylogenetics(46). The second option in C) is to choose a single best model, with the most obvious concern being overfitting, and parametric bootstrapping is not generally computationally feasible here. An accepted approach for choosing a

single best model is to use the exact order based method of (45) which finds the globally most probable posterior DAG.

All modeling results were carried out in R (47) using an R library called **abn** developed by the authors for the purpose of analyzing epidemiological data. This software is freely available for download from CRAN.

Results

Multivariable regression

Table 1 shows the variables in the optimal model from the stepwise multivariable regression analysis (see Web Appendix 6). There are 12 covariates, many of which have low p-values. To identify spurious covariates, arising from overfitting, parametric bootstrapping was used to generate 10,000 data sets from the optimal model. The parametric bootstrapping results provide convincing evidence that each of the 12 identified covariates is robust in terms of being statistically associated with the presence of diarrhea (Web Appendix 6). This regression model can be represented as a DAG where each of the explanatory variables is a node with an arc directed towards the node for the response variable (Figure 2).

Additive Bayesian Network

Three different, although related, sets of results are presented, all with the same goal of identifying those variables directly dependent with the presence of diarrhea.

Heuristic search across 13 variables

The standard local heuristic search(27) was applied to the subset of 12 covariates identified in the optimal multiple regression model. A (directed) majority consensus ABN model was constructed by pooling results across 20,000 heuristic searches and was sufficient for robust results (see Web Appendix 7). This summary network (see Web Appendix 7) identifies “Dry Pit Latrine” and access to an atypical - “Other Water Source”, as directly dependent with diarrhea. Figure 3 shows an undirected majority consensus ABN constructed from the same 20,000 heuristic searches and this now additionally has “No Formal Garbage Collection” as directly dependent with diarrhea, although its structural support is relatively weaker, in terms of how often it was chosen for inclusion each each locally optimal DAG (see Web Appendix 7), than the other two variables. Posterior density estimates for the three variables directly dependent with diarrhea can be found in Web Appendix 8. Note that the odds ratio estimates in Table 1 and posterior densities will be identical in any ABN which has only these three variables with arcs to diarrhea.

Exact search for most probable DAG across 13 variables

The most probable posterior DAG was identified using the exact method of (45), again on the reduced set of 13 variables (specific details are given Web Appendix 9). This exact search identifies a maximal ABN which has “Dry Pit Latrine” and access to an atypical - “Other Water Source”, as directly dependent with diarrhea but not “No Formal Garbage Collection”. The goodness of fit of this model is -105028.4 (log marginal likelihood) and it has 32 arcs in total. During the previous 20,000 search heuristics across DAGs a number of models with improved

goodness of fit were identified e.g. -105025.9 with 34 arcs, which demonstrates the bias toward parsimony in order based searches, as the fewer arcs a DAG has the more orders it will be consistent with.

Heuristic search across all 41 variables

The standard heuristic search (27) over 41 variables was not computationally feasible, and similarly for the exact order-based method. By necessity an ad-hoc approach was instead utilized by adding several constraints to the standard heuristic search (see Web Appendix 10 for details). A majority consensus ABN model was constructed by pooling results across $\approx 500,000$ separate searches. This model identifies the same three variables as directly dependent with diarrhea as in the undirected majority consensus network with 13 variables. The ABN model supports 182 inter-dependencies between the 41 variables, where 179 are dependencies indirectly related to the presence of diarrhea, that is between variables which can potentially affect disease presence, but only through their relationships with other variables (see Web Appendix 11 for a detailed description of the model).

Discussion

The objective of the analyzes presented was to identify potential determinants of the presence of diarrhea, and in particular to contrast results using standard multivariable regression with that of an epidemiological systems approach utilizing additive Bayesian networks.

Comparison of methods

Table 1 along with Figures 2 and 3 show clearly that the two approaches provide very different, though overlapping, results. The ABN results suggest that most - 9 out of 12 - of the covariates identified in the multivariable regression analyzes, whilst associated with the presence of diarrhea, are only indirectly rather than directly related with this outcome.

Additive Bayesian network models are simply multivariate extensions of standard multivariable regression - nothing more. The single key conceptual difference is that ABN models are multidimensional and consider all relationships between all variables simultaneously. It is therefore intuitively reasonable to expect both approaches - ABN and multivariable regression - to both identify in common those variables with the strongest degree of statistical support (and irrespective of whether different goodness of fit metrics or inferential framework are used e.g. AIC or marginal likelihood, Bayesian or non-Bayesian). This is exactly what was found in the analyzes presented (Table 1) - the three variables with lowest p-values in the multivariable regression using AIC are also those supported as directly dependent with diarrhea in the ABN results.

In the multivariable regression analyzes “Number Of Rooms” was identified as associated with diarrhea, and with a p-value sufficiently low (0.007) to be typically considered as strong statistical evidence, and this was further supported through the bootstrapping results. This variable has many direct dependencies in the ABN (Figure 3 and Web Appendix 11) but only with variables other than diarrhea. Biologically, “Number Of Rooms” cannot be directly dependent with diarrhea as while it is likely to be related to the living and environmental conditions in a

household - and the ABN model provides empirical evidence of this - intuitively, it cannot contribute directly to diarrhea infection. This suggests that “Number Of Rooms” has been identified in the multivariable regression model as a result of association induced with diarrhea through a network of inter-dependencies across the disease system. This highlights the difficulty of interpretation in the traditional multivariable regression model where it is possible for variables with low p-values to be identified as associated with disease but which are likely to be only indirectly related to the outcome variable.

Biological interpretation of results

The decreased risk of diarrhea from dry pit latrines suggests that infectious enteric pathogens are efficiently removed from faecal-oral transmission cycles. Using a univariate regression model, the absence of a toilet in a household was a substantially greater risk factor for diarrhea (OR 5.7, 95% CI 5.07, 6.50) than the presence of any type of toilet. The absence of arcs connecting “No Toilet In The Household” to childhood diarrhea suggests, however, that those houses lacking toilets have a network of confounding factors that modify the risk of enteric infection. For example, a fuller description of the disease system (Web Appendix 11) shows that the absence of a toilet is dependent with other descriptions of the household such as the absence of an electrical connection (as an indicator of socio-economic status), with certain water sources, especially those that do not require infrastructure, such as ponds, streams and springs, and also no formal garbage collection. What might be deduced from this is that the lack of a toilet

is an index of lower socio-economic status and therefore living conditions. Socio-economic status is itself associated with education levels and certain behaviors, e.g. unhygienic practices that increase the risk of diarrhea.

Using “other” sources of water is comparatively unusual in the data (2.1% of households). Given the breadth of options covered across the other five categories of water sources (see Web Appendix 2), “Other Water Source” is somewhat vague, but includes buying in water from a local seller and collecting water from a tanker-supplied public standpipe. Water that is further removed from either natural water courses or pipe/well systems may be less susceptible to leakage between sewerage systems and the water table and thereby at risk of contamination with water-borne pathogens. Examination of the 41 variable ABN model shows connections from each of the water sources to at least the alternative sewerage connections or to the type of toilet, demonstrating the tight interlinkage between these three risk categories in determining the exposure of children to enteric pathogens. The relatively large uncertainty in the log odds estimate between “other” water sources and diarrhea (Web Appendix 8) is likely due to both the rarity of this water source in the data and also the ambiguous definition of this variable. The Web Appendix 3 contains a list of additional references and a summary of previous variables associated with diarrhea, including age which is also briefly discussed.

In the language of Hernan et al (14) the components of a disease system are the disease outcome, an exposure (that directly results in disease) and a series of confounding variables that act either on the exposure or both exposure and disease outcome. The use of survey data, as in the PSLSM, does not necessarily

contain the sort of proximal exposure that results in infection (contact with enteric pathogens), so much as a collection of variables that might be considered common to both exposure (infection) and disease (pathogenesis). It is, therefore, arguably less useful to discuss causality in respect of such an exploratory model, however, what is both possible and useful is the partitioning of factors into those that are directly or indirectly dependent with disease outcome. To illustrate this point, consider the type of toilet that is present in a house: there are several types of toilet in the PSLSM data, but these are all variants of the same underlying theme of how the household disposes of human excreta. It is not practical to question whether a flushing toilet is casual of diarrheal disease so much as the relative risk of disease given the different types of toilet present - the actual exposure is still the contact with enteric pathogens, which can then be stratified (assuming such data exists) by type of toilet. In this context “toilet” is a confounder (according to Hernan et al) that can be used to stratify the more proximal exposure.

Following through a series of DAGs as subsets of the more complete system may offer a closer parallel to causal models that are based on expert opinion (Web Appendix 12). Assuming that the absence of formal garbage collection is the exposure leading to diarrheal disease, four alternative routes involving four additional confounding factors were compared. The reason for selecting “no formal garbage collection” was because this variable was identified as dependent with diarrhea in all but the most probable DAG. There was little change in the odds of diarrhea when no formal garbage collection was combined with different combinations of other confounders in the alternative DAGs. The interesting exception is the combination of “other” sources of water and no formal garbage

collection. With the more defined water sources, the odds of diarrhea in houses lacking formal garbage collection is approximately 1.2, however this rises to 1.3 when water sources are “other”. The explanation for this may be the results of lack of running water to remove the build-up of refuse (and potentially excreta) when garbage is not regularly removed. The accumulation of refuse that has no formal disposal, and its removal is presumably irregular, provides a permissive breeding ground for enteric pathogens. This hypothesis requires more targetted studies, however, as neither garbage nor water are of themselves the “cause” of diarrheal disease despite both being likely sources of enteric pathogens.

Introducing prior causal knowledge

Bayesian network modeling is typically concerned with automated structure discovery - searching for a DAG which best describes the statistical relationships in observational data, in this case the PSLSM. In causal inference, on the other hand, the focus is typically on testing whether a given set of assumptions is sufficient for quantifying causal effects from observational data, conditional on a causal diagram which encodes all the relevant domain-specific assumption(12, 14). The former is objective - empirically derived DAGs - but lacks a causal component, while the latter provides causal insight but whose weakness is the potentially subjective justification of the causal diagram. An obvious question is, therefore, how can prior causal knowledge be integrated into automated structure discovery.

A very rudimentary approach which repeats the previous heuristic ABN analyzes (across 13 variables) by introducing some simple common sense prior causal constraints is given in Web Appendix 13. This modeling prohibits arcs ema-

nating from the diarrhea node and prevents “Number Of Rooms” being directly dependent with diarrhea. This simple informative structural prior now gives an undirected majority consensus DAG with the same three variables as previously identified as directly dependent with diarrhea, but where support for an arc connecting “No Formal Garbage Collection” to diarrhea is now 100% (appears in all 20,000 searches) whereas using the previous uninformative structural prior this was only 58% and it also didn’t appear in the previous directed majority consensus network (Web Appendix 7). The use of directional constraints alters the model search space - as the data can discriminate between DAGs where the arcs in these constraints are reversed provided these are in different equivalence classes (see Web Appendix 1) - and so may provide different results. The key question is whether imposing such directional constraints/informative prior - motivated by causal considerations - is conceptually reasonable. An open question.

An alternative approach - and one which seems preferable given the complications of likelihood equivalence - is suggested in (27) (p.224). Rather than use an informative structural prior it is proposed to append onto the observed data additional - and likely highly incomplete - synthetic observations which reflect causal beliefs. The structure learning process is then applied to all of the data as usual, except with the additional functionality necessary to marginalize over missing data (e.g. (48)). This is an elegant approach but its feasibility and practicality in respect of ABN modeling is an open question and an exciting area of future work.

Acknowledgments

Author affiliations: Section of Epidemiology, University of Zurich, Zurich, Switzerland (Fraser I. Lewis); Fogarty International Center, National Institutes of Health, Bethesda, USA (Benjamin J. J. McCormick).

BJJM was supported by the Bill and Melinda Gates Foundation and the Fogarty International Center of the National Institutes of Health as part of the Mal-ED network.

The authors thank Safdar Parvez, ADB Manila, for assisting in gaining access to the PSLSM data.

Conflict of interest: none declared.

References

1. Black RE, Cousens S, Johnson HL, et al. Global, regional, and national causes of child mortality in 2008: a systematic analysis. *Lancet*. 2010; 375(9730):1969–1987.
2. Bryce J, Boschi-Pinto C, Shibuya K, et al. WHO estimates of the causes of death in children. *Lancet*. 2005;365(9465):1147–1152.
3. Ahiadeke C. Breast-feeding, diarrhoea and sanitation as components of infant and child health: A study of large scale survey data from Ghana and Nigeria. *J Biosoc Sci*. 2000;32(1):47–61.
4. de Souza A, Peterson K, Cufino E, et al. Relationship between health services, socioeconomic variables and inadequate weight gain among Brazilian children. *B World Health Organ*. 1999;77(11):895–905.

5. Hatt L, Waters H. Determinants of child morbidity in Latin America: A pooled analysis of interactions between parental education and economic status. *Social Science & Medicine*. 2006;62(2):375–386.
6. Jones L, Griffiths P, Adair L, et al. A comparison of the socio-economic determinants of growth retardation in South African and Filipino infants. *Public Health Nutr*. 2008;11(12):1220–1228.
7. Boerma JT, Black RE, Sommerfelt AE, et al. Accuracy and completeness of mothers recall of diarrhea occurrence in preschool-children in demographic and health surveys. *Int J Epidemiol*. 1991;20(4):1073–1080.
8. Galea S, Riddle M, Kaplan GA. Causal thinking and complex system approaches in epidemiology. *Int J Epidemiol*. 2010;39(1):97–106.
9. Fenner L, Egger M, Gagneux S. Annie darwin’s death, the evolution of tuberculosis and the need for systems epidemiology. *Int J Epidemiol*. 2009; 38(6):1425–1428.
10. Lusi AJ, Attie AD, Reue K. Metabolic syndrome: from epidemiology to systems biology. *Nat Rev Genet*. 2008;9(11):819–830.
11. Roux AVD. From the American college of epidemiology annual meeting 2006 - integrating social and biologic factors in health research: A systems view. *Ann Epidemiol*. 2007;17(7):569–574.
12. Pearl J. Causal diagrams for empirical research. *Biometrika*. 1995;82(4):669–688.

13. Hernan MA, Robins JM. Estimating causal effects from epidemiological data. *J Epidemiol Commun H.* 2006;60(7):578–586.
14. Hernan MA, Hernandez-Diaz S, Werler MM, et al. Causal knowledge as a prerequisite for confounding evaluation: An application to birth defects epidemiology. *Am J Epidemiol.* 2002;155(2):176–184.
15. Eisenberg JNS, Desai MA, Levy K, et al. Environmental determinants of infectious disease: A framework for tracking causal links and guiding public health research. *Environ Health Persp.* 2007;115(8):1216–1223.
16. Howard G, McClure LA, Moy CS, et al. Imputation of incident events in longitudinal cohort studies. *Am J Epidemiol.* 2011;174(6):718–726.
17. Johnson HL, Liu L, Fischer-Walker C, et al. Estimating the distribution of causes of death among children age 1-59 months in high-mortality countries with incomplete death certification. *Int J Epidemiol.* 2010;39(4):1103–1114.
18. Shultz A, Omollo J, Burke H, et al. Cholera outbreak in Kenyan refugee camp: Risk factors for illness and importance of sanitation. *Am J Trop Med Hyg.* 2009;80(4):640–645.
19. Phillips G, Lopman B, Rodrigues L, et al. Asymptomatic rotavirus infections in england: Prevalence, characteristics, and risk factors. *Am J Epidemiol.* 2010;171(9):1023–1030.
20. Breiman L. Statistical modeling: The two cultures. *Stat Sci.* 2001;16(3):199–215.

21. Cox DR, Efron B, Hoadley B, et al. Statistical modeling: The two cultures - comments and rejoinders. *Stat Sci.* 2001;16(3):216–231.
22. Arbogast PG, Ray WA. Performance of disease risk scores, propensity scores, and traditional multivariable outcome regression in the presence of multiple confounders. *Am J Epidemiol.* 2011;174(5):613–620.
23. Babyak MA. What you see may not be what you get: A brief, nontechnical introduction to overfitting in regression-type models. *Psychosom Med.* 2004; 66(3):411–421.
24. Raftery AE, Madigan D, Hoeting JA. Bayesian model averaging for linear regression models. *J Am Stat Assoc.* 1997;92(437):179–191.
25. Posada D, Buckley TR. Model selection and model averaging in phylogenetics: Advantages of akaike information criterion and Bayesian approaches over likelihood ratio tests. *Systematic Biol.* 2004;53(5):793–808.
26. Friedman N, Goldszmidt M, Wyner A. Data analysis with Bayesian networks: A bootstrap approach. In: *Proc. Fifteenth Conference on Uncertainty in Artificial Intelligence. San Francisco: Morgan Kaufmann.* 206–215.
27. Heckerman D, Geiger D, Chickering DM. Learning Bayesian networks - the combination of knowledge and statistical-data. *Mach Learn.* 1995;20(3):197–243.
28. Needham CJ, Bradford JR, Bulpitt AJ, et al. A primer on learning in Bayesian networks for computational biology. *PLoS Comput Biol.* 2007; 3(8):e129.

29. Djebbari A, Quackenbush J. Seeded Bayesian networks: Constructing genetic networks from microarray data. *BMC Syst Biol.* 2008;2–57.
30. Poon AFY, Lewis FI, Pond SLK, et al. Evolutionary interactions between n-linked glycosylation sites in the HIV-1 envelope. *PLoS Comput Biol.* 2007; 3(1):110–119.
31. Poon AFY, Lewis FI, Pond SLK, et al. An evolutionary-network model reveals stratified interactions in the V3 loop of the HIV-1 envelope. *PLoS Comput Biol.* 2007;3(11):2279–2290.
32. Poon AFY, Lewis FI, Frost SDW, et al. Spidermonkey: rapid detection of co-evolving sites using Bayesian graphical models. *Bioinformatics.* 2008; 24(17):1949–1950.
33. Lycett SJ, Ward MJ, Lewis FI, et al. Detection of mammalian virulence determinants in highly pathogenic avian influenza H5N1 viruses: Multivariate analysis of published data. *J Virol.* 2009;83(19):9901–9910.
34. Jansen R, Yu HY, Greenbaum D, et al. A Bayesian networks approach for predicting protein-protein interactions from genomic data. *Science.* 2003; 302(5644):449–453.
35. Hodges AP, Dai D, Xiang Z, et al. Bayesian network expansion identifies new ROS and Biofilm regulators. *PLoS ONE.* 2010;5(3):e9513.
36. Dojer N, Gambin A, Mizera A, et al. Applying dynamic Bayesian networks to perturbed gene expression data. *BMC Bioinformatics.* 2006;7:249.

37. Lewis FI, Brulisauer F, Gunn GJ. Structure discovery in Bayesian networks: An analytical tool for analysing complex animal health data. *Prev Vet Med.* 2011;100(2):109–115.
38. Cooper GF, Herskovits E. A Bayesian method for the induction of probabilistic networks from data. *Mach Learn.* 1992;9(4):309–347.
39. Milns I, Beale CM, Smith VA. Revealing ecological networks using Bayesian network inference algorithms. *Ecology.* 2010;91(7):1892–1899.
40. Kim SYst, Imoto S, Miyano S. Inferring gene networks from time series microarray data using dynamic bayesian networks. *Brief Bioinform.* 2003;4(3):228–235.
41. Rijmen F. Bayesian networks with a logistic regression model for the conditional probabilities. *Int J Approx Reason.* 2008;48(2):659–666.
42. Boettcher SG, Dethlefsen C. deal: A package for learning Bayesian networks. *J Stat Softw.* 2003;8(20):1–40.
43. Smith AFM. Bayesian computational methods. *Philos T Roy Soc A.* 1991;337(1647):369–386.
44. Friedman N, Koller D. Being Bayesian about network structure. a bayesian approach to structure discovery in bayesian networks. *Mach Learn.* 2003;50(1-2):95–125.
45. Koivisto M, Sood K. Exact Bayesian structure discovery in Bayesian networks. *J Mach Learn Res.* 2004;5:549–573.

46. Ronquist F, Huelsenbeck JP. Mrbayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics*. 2003;19(12):1572–1574.
47. R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2006. ISBN 3-900051-07-0.
48. Spiegelhalter DJ, Lauritzen SL. Sequential updating of conditional probabilities on directed graphical structures. *Networks*. 1990;20(5):579–605.

Table 1: Results of stepwise regression search and ABN analyzes (with all three variables directly dependent with diarrhea). OR - odds ratio. P-values are from Type III Chi-squared tests. The ORs are marginal, for example for Dry Pit Latrine the OR = 0.66 which is relative to not having a Dry Pit Latrine (ignoring all other covariates), for the continuous variables (Age and Number Of Rooms) the ORs are in respect of a one unit increase.

Indicator	OR	p-value	Bayesian OR
Child Age	1.05	0.012	-
Number Of Rooms	0.95	0.007	-
Sex			
- Male	1.14	0.012	-
Dwelling Type			
- Part Of Compound	0.70	0.020	-
Source Drinking Water			
- Piped	0.85	0.011	-
- Canal/River/Stream	0.65	0.0087	-
- Spring	0.76	0.13	-
- Other	0.46	0.0016	0.49
Type Of Toilet			
- Flush, Connected To Open Drain	0.84	0.046	-
- Dry Pit Latrine	0.66	<0.0001	0.67
Connection To Sewerage			
- Yes, Covered Drains	0.72	0.064	-
Organizer Of Garbage Collection From House			
- No Formal System	1.23	0.0048	1.32

Figure 1: Example of an additive Bayesian network model comprising of three binary random variables (X_1, X_2, X_5) and two continuous (Gaussian) variables (X_3, X_4). The model for each node is a generalized linear regression with identity or logit link function as appropriate. Let π_i for $i = 1, 2, 5$ denote the probability of observing a success: $P(X_i = 1) = 1 - P(X_i = 0)$, and μ_i the mean of random variable X_i for $i = 3, 4$. X_2 is independent of the other variables with $\log\{\pi_2/(1 - \pi_2)\} = \beta_{2,0}$; X_4 is conditionally dependent upon X_2 with $\mu_4 = \beta_{4,0} + \beta_{4,1}X_2$; X_5 is conditionally dependent upon X_4 with $\log\{\pi_5/(1 - \pi_5)\} = \beta_{5,0} + \beta_{5,1}X_4$; X_3 is conditionally dependent upon X_4 with $\mu_3 = \beta_{3,0} + \beta_{3,1}X_4$; and X_1 is conditionally dependent upon X_2, X_3, X_5 with $\log\{\pi_1/(1 - \pi_1)\} = \beta_{1,0} + \beta_{1,1}X_2 + \beta_{1,2}X_3 + \beta_{1,3}X_5$.

Figure 2: Final model in stepwise (forwards and backwards) multivariable regression search depicted as a DAG. Parametric bootstrapping statistically supports all 12 covariates in this model. Ovals are continuous variables, squares discrete.

Figure 3: Undirected majority consensus ABN model constructed by pooling results across 20,000 heuristic searches. Only three variables: “No Formal Garbage Collection”, access to a “Dry Pit Latrine” and access to an atypical - “Other Water Source” are supported as directly dependent with the presence of diarrhea.